

A deep learning model to predict RNA-seq expression of tumours from whole slide images

Benoît Schmauch¹, Alberto Romagnoni^{1,^}, Elodie Pronier^{1,^}, Charlie Saillard¹, Pascale Maillé², Julien Calderaro², Aurélie Kamoun¹, Meriem Sefta¹, Sylvain Toldo¹, Mikhail Zaslavskiy¹, Thomas Clozel¹, Matahi Moarii¹, Pierre Courtiol^{1,‡}, Gilles Wainrib^{1,‡}

¹: Owkin Lab, Owkin, Inc. New York, NY USA.

²: INSERM U955, Team "Pathophysiology and Therapy of Chronic Viral Hepatitis and Related Cancers", Crêteil, France; APHP, Department of Pathology, Hôpital Henri Mondor, Université Paris-Est, Crêteil, France.

^{^,‡}: These authors contributed equally.

Corresponding authors:

Gilles Wainrib:

Owkin Lab, Owkin, Inc. New York, NY, 10003 USA.

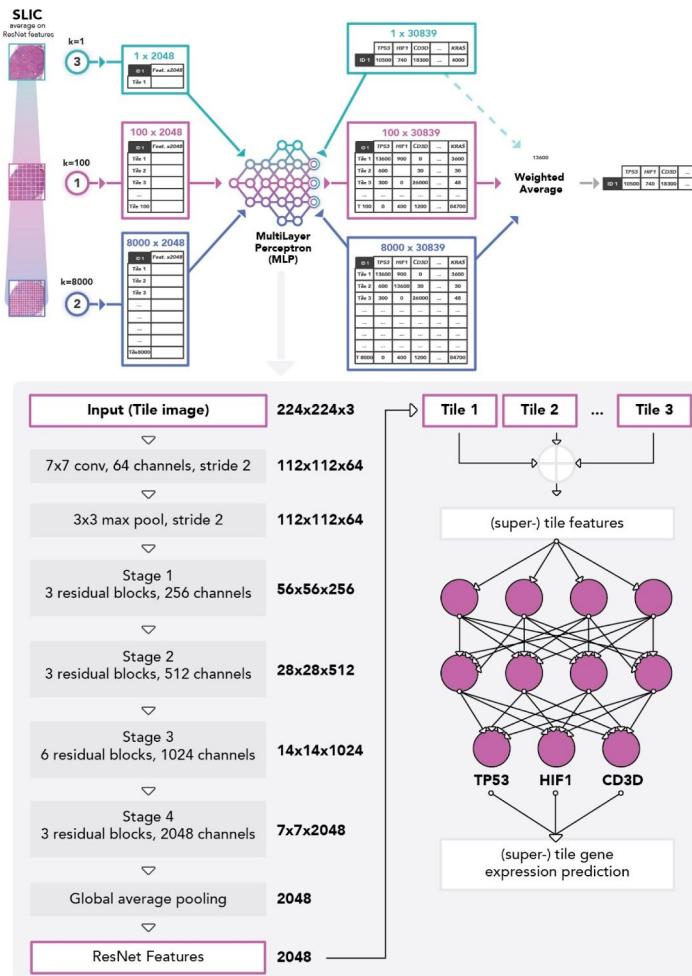
gilles.wainrib@owkin.com

Benoît Schmauch:

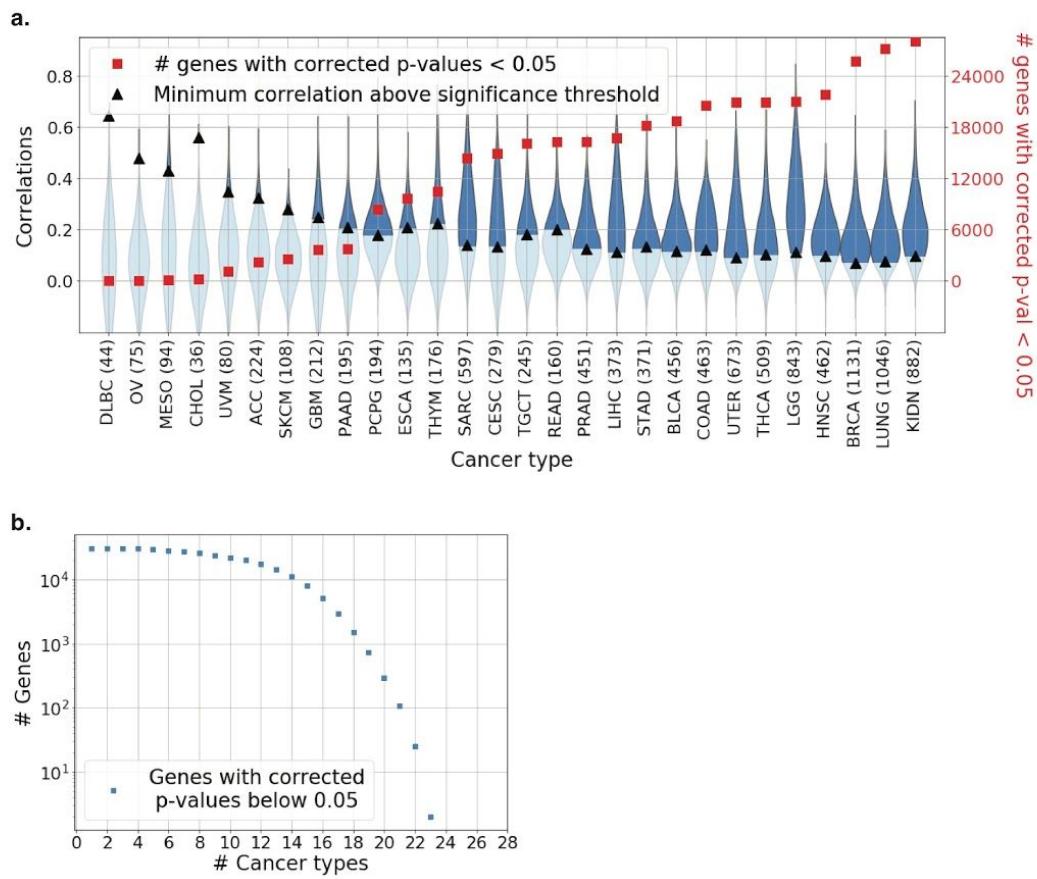
Owkin Lab, Owkin, Inc. New York, NY, 10003 USA.

benoit.schmauch@owkin.com

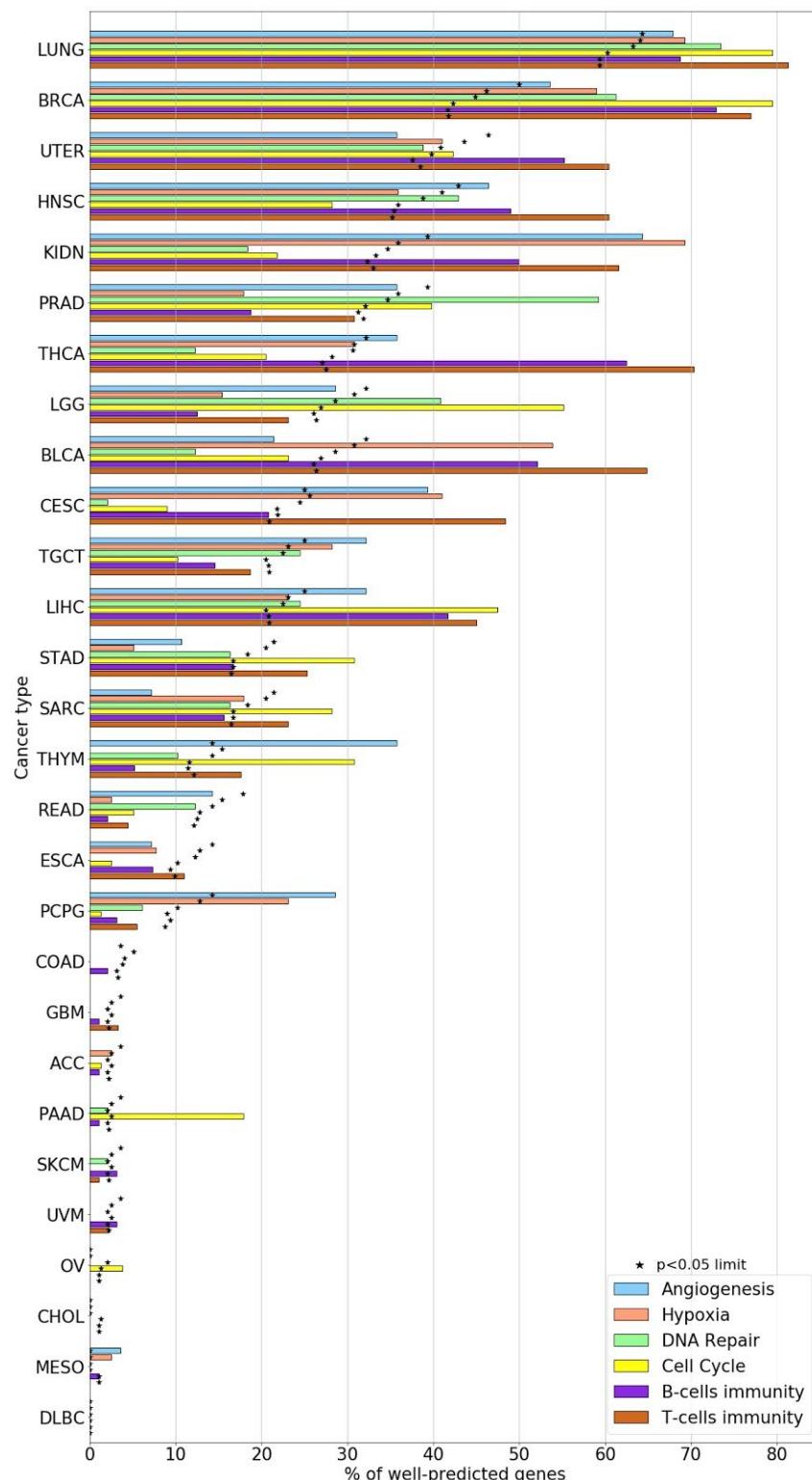
Supplementary Information



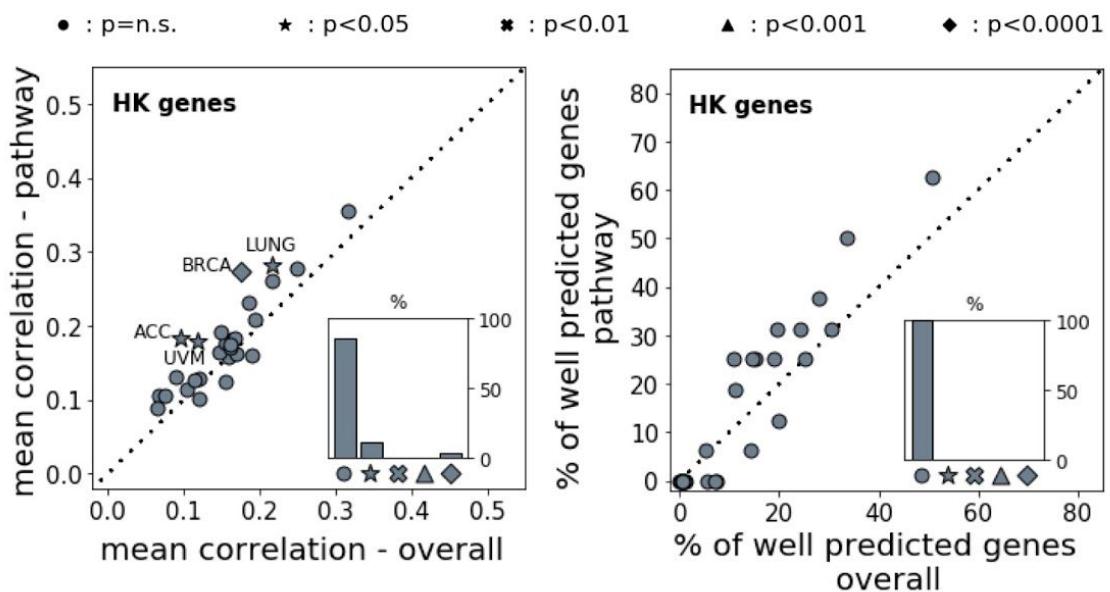
Supplementary Figure 1. Preprocessing and model structure. The whole-slide image is divided into $112 \times 112 \mu\text{m}$ tiles (224×224 pixels). A maximum of 8,000 such tiles (excluding those covering the background) are kept. A 50-layer ResNet, pretrained on the ImageNet dataset, is used to extract 2,048 features from each tile image (bottom left). A preprocessing algorithm, inspired by simple linear iterative clustering (SLIC) is applied to the 8,000 tiles to produce k super-tiles. An average is obtained for each super-tile, at the level of the 2,048 ResNet features. The number of clusters is decided according to the task (color and numerical code as in Fig. 1). A multilayer perceptron (bottom right) is applied to each cluster of the slide. The last layer of the model encodes the transcriptomic representation described in the text. This representation is then used to produce a prediction per cluster and per Ensembl gene of the RNA-Seq dataset. Finally, a weighted average (described in the Methods) provides the output prediction of gene expression associated with the slide.



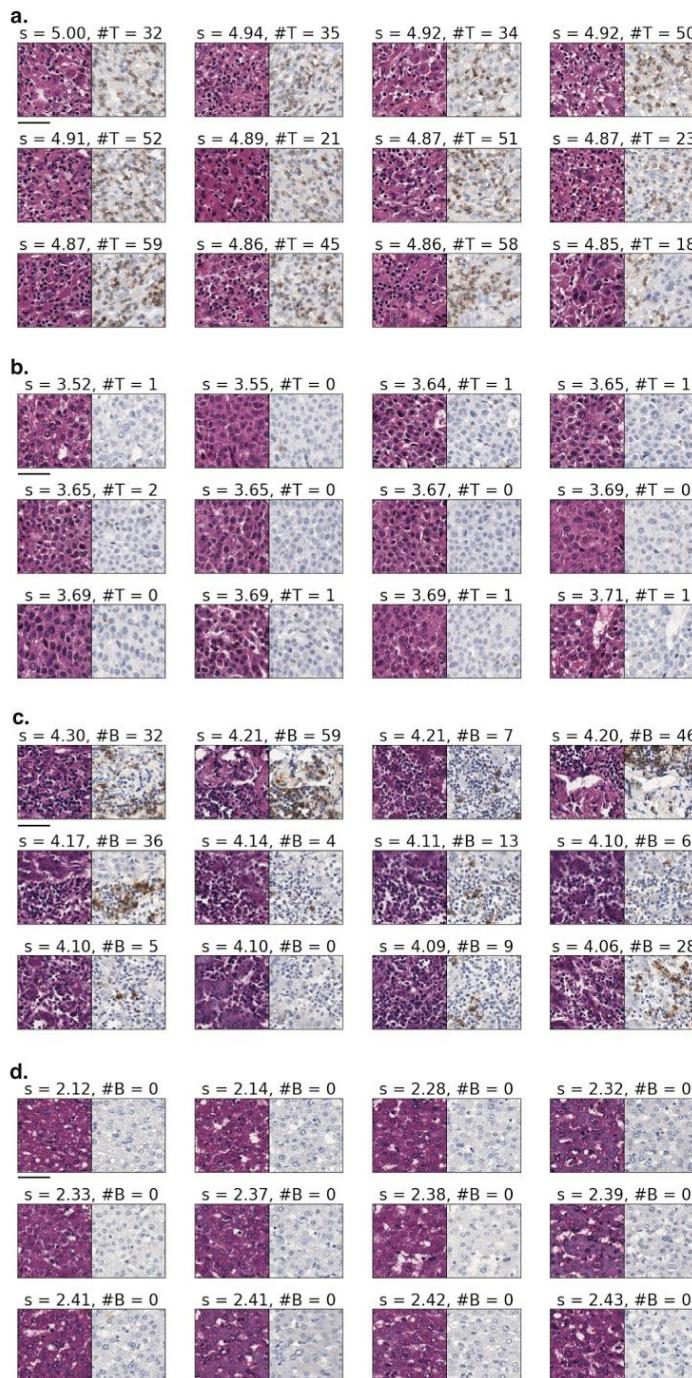
Supplementary Figure 2. Predictions of gene expression after Benjamini–Hochberg correction for the testing of multiple hypotheses. a. Distribution of Pearson's correlation coefficients R (left axis, blue violin plots) and the number of coding and non-coding genes (right axis, red squares) using Benjamini–Hochberg-corrected p-values < 0.05 (one-sided empirical p-value, as described in Methods section), for twenty eight cancer types from the TCGA. Black triangles indicate the minimal correlation coefficient required for significance in any given dataset. **b.** Number of coding and non-coding genes for which expression was significantly well-predicted for a given number of cancer types, as a function of the number of cancers.



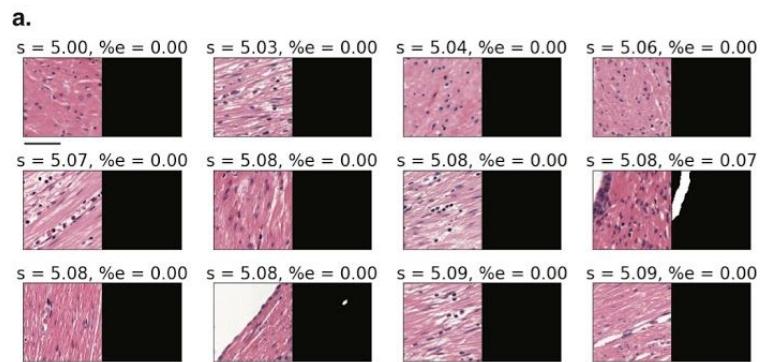
Supplementary Figure 3. Percentage of genes for which expression was well-predicted in six common pathways of carcinogenesis. Percentage of genes for which expression was well-predicted (as in Fig.2) for the six studied hallmark pathways for cancer. The black stars indicate the percentage required for each cancer dataset and each pathway to be considered significantly better predicted than a corresponding random list of genes of the same length as the pathway gene list (see Methods section).



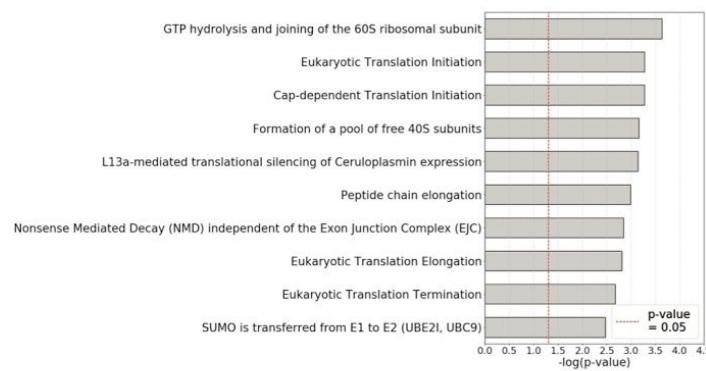
Supplementary Figure 4. Prediction of a housekeeping genes signature. As in Fig.3, for the housekeeping gene signature described in Supplementary Table 4. *Left panel.* The indicated statistical significance refers to the probability of obtaining a correlation $R > R_p$ in the distribution of correlations for random lists, for each given cancer type. Insets show the percentages of the different cases of statistical significance between cancer types. *Right panel.* As in the left panel, but for the percentage of genes for which expression was of well-predicted (as defined in the text and in Fig. 2). HK = housekeeping.



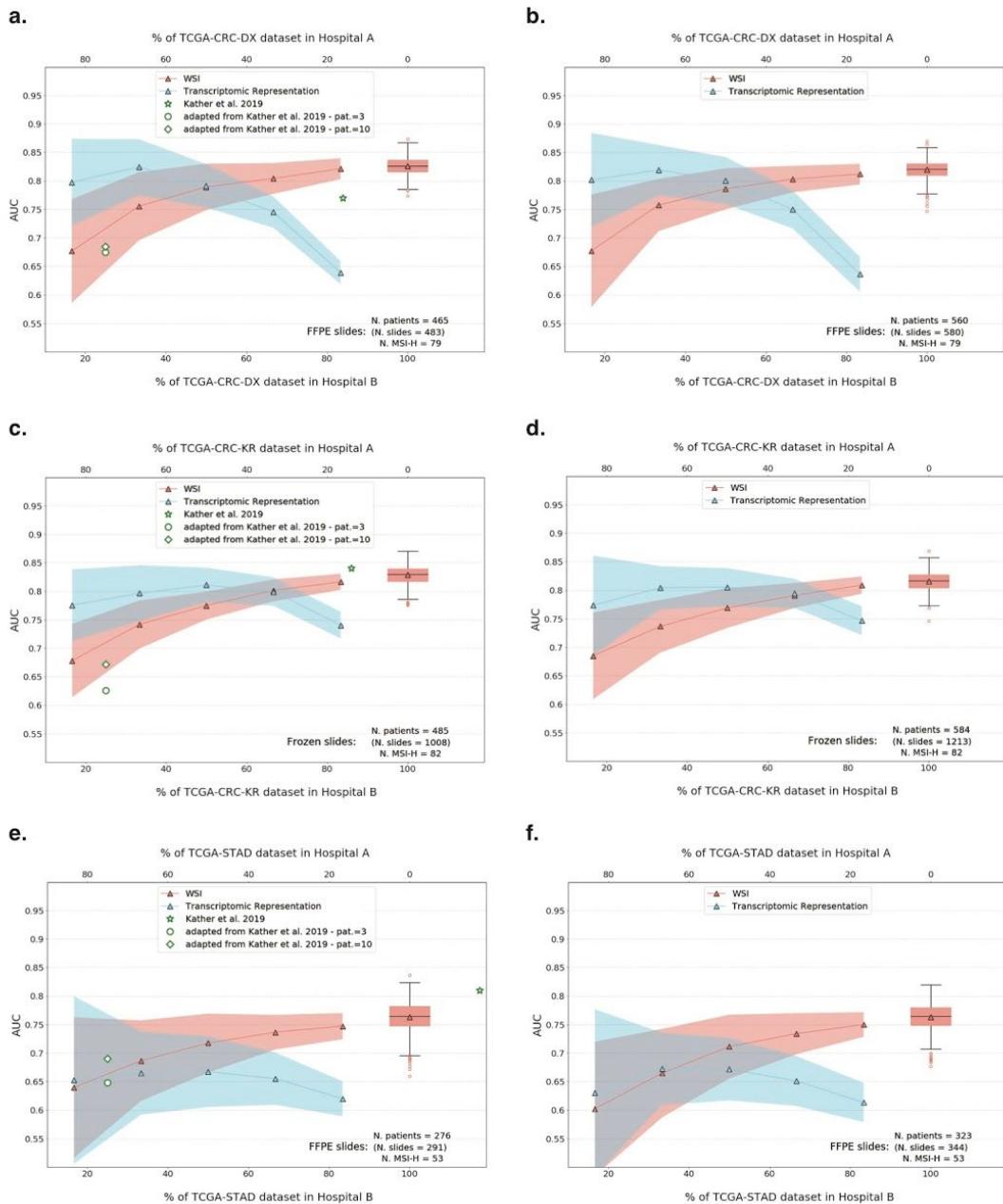
Supplementary Figure 5. Most predictive tiles for CD3-receptor encoding genes. a. Extraction of the tiles associated with the highest score for CD3⁺-encoding genes expression predicted from the HE/CD3 slide. S = Score corresponding to the log expression score of each tile; #T= Number of T cells per tile as determined with QuPath from the CD3 IHC. **b.** As in **a** but for the tiles with the lowest score for CD3 expression. **c.** Extraction of the tiles associated with the highest score for CD19 and CD20 expression predicted from the HE/CD20 slide. S = Score corresponding to the log expression score of each tile; #B= Number of B cells per tile as determined with QuPath from the CD20 IHC. **d.** As in **b** but for the tiles with the lowest score for CD3 expression. Scale bar: 100μm. Representative tiles from one double-stained slide (n=28,123 tiles)



Supplementary Figure 6. Least predictive tiles for epithelium associated genes.
Extraction of the tiles associated with the lowest score for epithelium-associated genes expression from the same slide as in Fig. 5. S = Score corresponding to the log expression score of each tile; %e = fraction of pixels marked as belonging to the epithelium. Scale bar: 100 μ m. (n= 21,714 tiles).

a.

Supplementary Figure 7: Computational pathway analysis of the best-predicted genes among the microsatellite stable patients of the TCGA-CRC cohort.



Supplementary Figure 8. Prediction of microsatellite instability status using transfer learning from transcriptomic representation. Change in area under the ROC curve (AUC) for the model based on the transcriptomic representation learned in Hospital A and trained on Hospital B (blue) and the model directly based on WSI images from Hospital B (red), as a function of the fraction of the dataset used in the two hospitals, for the different datasets used in the present study, with or without MSI-Low (MSI-L) patients included in the negative class together with MSS patients ($n = 50$ different data splits between the hospitals, averaged over 10 different 3-folds CV; solid lines and triangles: mean over splits; shaded area: 68% confidence interval). Whenever comparison is possible, we also display the result reported by Kather et al.⁵⁵ and that obtained from adapting the same method with 25% of the data in Hospital B (see Methods section for further details). **a.** and **b.** TCGA-CRC-DX, respectively without and with MSI-L patients.

c. and **d.** TCGA-CRC-KR, resp. without and with MSI-L patients. **e.** and **f.** TCGA-STAD, resp. without and with MSI-L patients. Boxplot in each panel represents the distribution of AUC values (500 3-folds CV) over the whole dataset, for the model based on WSI images (box: interquartile range (IQR); horizontal line: median; whiskers: 1.5 times IQR, triangle: mean; open circles: outliers).

Project ID	Disease Name	# Patients	# Samples
BRCA	Invasive breast carcinoma	1057	1131
LUNG	Lung adenocarcinoma (LUAD)	944	1046
	Lung squamous cell carcinoma (LUSC)		
KIDN	Kidney chromophobe carcinoma (KICH)	843	882
	Kidney renal clear cell carcinoma (KIRC)		
	Kidney renal papillary cell carcinoma (KIRP)		
LGG	Brain lower-grade glioma	485	843
SARC	Sarcoma	252	597
UTER	Uterine carcinosarcoma (UCS)	558	673
	Uterine corpus endometrial carcinoma (UCEC)		
THCA	Thyroid carcinoma	497	509
COAD	Colon adenocarcinoma	445	463
HNSC	Head and neck squamous cell carcinoma	423	462
BLCA	Bladder urothelial carcinoma	383	456
PRAD	Prostate adenocarcinoma	399	451
LIHC	Liver hepatocellular carcinoma	359	373
STAD	Stomach adenocarcinoma	350	371
CESC	Cervical squamous cell carcinoma and endocervical adenocarcinoma	267	279
TGCT	Testicular germ cell tumors	149	245
ACC	Adrenocortical carcinoma	55	224

GBM	Glioblastoma multiforme	96	212
PAAD	Pancreatic adenocarcinoma	175	195
PCPG	Pheochromocytoma and paraganglioma	175	194
THYM	Thymoma	116	176
READ	Rectum adenocarcinoma	159	160
ESCA	Esophageal carcinoma	133	135
SKCM	Skin, cutaneous melanoma	97	108
MESO	Mesothelioma	74	94
UVM	Uveal melanoma	80	80
OV	Ovarian serous cystadenocarcinoma	74	75
DLBC	Lymphoid neoplasm diffuse large B cell lymphoma	44	44
CHOL	Cholangiocarcinoma	36	36

Supplementary Table 1: TCGA dataset, detailed information. A matched WSI RNA-Seq data pair is considered here to be a sample.

Rank of tile prediction	$n_{\text{tiles}} = 100$	$n_{\text{tiles}} = 8,000$
1	0.269	0.0210
1 to 2	0.126	
2 to 5	0.054	
5 to 10	0.026	
10 to 20	0.011	0.0099
20 to 50	0.004	0.0043
50 to 100	0.001	0.0021
100 to 200	-	0.0010
200 to 500	-	0.0004
500 to 1,000	-	0.0002
1,000 to 2,000	-	8×10^{-5}
2,000 to 5,000	-	2×10^{-5}

Supplementary Table 2: Coefficients of tile predictions in the weighted sum.

These coefficients are defining the slide-level gene expression predicted by the model during inference, as a function of their rank.

Angiogenesis	HALLMARK_ANGIOGENESIS, BIOCARTA_VEGF_PATHWAY, KEGG_VEGF_SIGNALING_PATHWAY, GO_ANGIOGENESIS
Hypoxia	HALLMARK_HYPOXIA, BIOCARTA_HIF_PATHWAY, GO_REGULATION_OF_CELLULAR_RESPONSE_TO_HYPOXIA
DNA repair	HALLMARK_DNA_REPAIR, REACTOME_DNA_REPAIR, GO_DNA_REPAIR
Cell cycle	BIOCARTA_CELLCYCLE_PATHWAY, KEGG_CELL_CYCLE, REACTOME_CELL_CYCLE, GO_CELL_CYCLE
B cell-mediated immunity	GO_B_CELL_MEDIATED_IMMUNITY
T cell-mediated immunity	REACTOME_ADAPTATIVE_IMMUNE_SYSTEM GO_ADAPTATIVE_IMMUNE_RESPONSE GO_REGULATION_OF_ADAPTATIVE_IMMUNE_RESPONSE

Supplementary Table 3: List of signatures from GSEA combined to define the six cancer pathways.

Ingenuity Canonical Pathways	-log(p-value)	Ratio	Genes
Th1 and Th2 Activation Pathway	14.1	0.0909	<i>CD247, CCR1, IL2RG, CD3E, IL12RB1, HAVCR2, CXCR3, CD8A, CD3D, CD3G, PIK3CG, IL10RA, CD86, IL2RA, VAV1, HLA-DPB1, HLA-DPA1</i>
iCOS-iCOSL Signaling in T Helper Cells	12.9	0.112	<i>PTPRC, CD247, CD3G, IL2RG, LCK, CD3E, PIK3CG, ZAP70, TRAT1, CD86, IL2RA, VAV1, CD3D, ITK</i>
Th1 Pathway	12.4	0.102	<i>CD247, CD3E, IL12RB1, HAVCR2, CXCR3, CD8A, CD3D, CD3G, PIK3CG, IL10RA, CD86, VAV1, HLA-DPB1, HLA-DPA1</i>
T Cell Receptor Signaling	10.6	0.1	<i>PTPRC, CD247, CD3G, LCK, PTPN7, CD3E, PIK3CG, ZAP70, VAV1, CD8A, CD3D, ITK</i>
Th2 Pathway	10.5	0.0855	<i>CCR1, CD247, CD3G, IL2RG, CD3E, IL12RB1, PIK3CG, CD86, IL2RA, VAV1, HLA-DPB1, CD3D, HLA-DPA1</i>
CD28 Signaling in T Helper Cells	10	0.0896	<i>PTPRC, CD247, CD3G, LCK, CD3E, WAS, PIK3CG, ZAP70, CD86, VAV1, CD3D, ITK</i>
Primary Immunodeficiency Signaling	8.87	0.16	<i>PTPRC, IL2RG, LCK, CD3E, ZAP70, CIITA, CD8A, CD3D</i>
CTLA4 Signaling in Cytotoxic T Lymphocytes	8.84	0.099	<i>CD247, CD3G, LCK, CD3E, PIK3CG, ZAP70, TRAT1, CD86, CD8A, CD3D</i>
Pathogenesis of Multiple Sclerosis	8.82	0.556	<i>CCR1, CCL4, CXCR3, CCL5, CXCL9</i>
Natural Killer Cell Signaling	7.68	0.0752	<i>CD247, LCK, SH2D1A, LAIR1, TYROBP, PIK3CG, ZAP70, VAV1, HCST, FCGR3A/FCGR3B</i>

Supplementary Table 4. Pathways well-predicted in at least 12 cancer types.

Canonical pathways with the best overlap with the 156 genes for which expression was well-predicted in at least 12 cancer types. P-values were calculated using Right-Tailed Fisher's Exact Test.

Ingenuity Canonical Pathways	-log(p-value)	Ratio	Genes
Cell Cycle Control of Chromosomal Replication	10.6	0.286	<i>MCM6, CDC45, CDT1, CDK16, CDC6, ORC6, CDC7, CDK1, MCM4, MCM3, MCM2, TOP2A, PRIM2, DBF4, ORC1, MCM7</i>
Mitotic Roles of Polo-Like Kinase	8.41	0.227	<i>KIF23, CDC20, PTTG1, PRC1, CDC7, CCNB2, PLK1, CDK1, CCNB1, PLK4, TGFB1, FBXO5, PKMYT1, KIF11, CDC25A</i>
Hepatic Fibrosis / Hepatic Stellate Cell Activation	8.39	0.134	<i>COL8A2, CCR5, COL10A1, COL4A2, COL1A2, COL5A1, COL16A1, TIMP1, TGFB1, PDGFRA, TIMP2, CXCL8, COL6A2, FGFR2, MMP2, COL1A1, COL6A3, TGFB3, TGFA, IL10RA, EDNRA, COL11A1, COL9A2, MMP9, COL3A1</i>
Th1 and Th2 Activation Pathway	7.68	0.128	<i>CCR5, IL2RG, HLA-DOA, CD3E, IL12RB1, IKZF1, HAVCR2, PIK3R5, HLA-DQA1, LGALS9, FGFR2, SPI1, CD3G, IL18, TGFB1, HLA-DMB, IL10RA, CD86, IL2RA, VAV1, JAG1, JAK3, NOTCH1, HLA-DPA1</i>
GP6 Signaling Pathway	6.85	0.141	<i>COL8A2, COL6A2, PIK3R5, COL10A1, FGFR2, COL4A2, LAMC2, COL16A1, COL5A1, COL1A2, COL1A1, COL6A3, SYK, LAMB1, FCER1G, COL11A1, COL9A2, LCP2, COL3A1</i>
Th2 Pathway	6.69	0.132	<i>CCR5, IL2RG, HLA-DOA, CD3E, IL12RB1, IKZF1, PIK3R5, HLA-DQA1, FGFR2, SPI1, CD3G, TGFB1, HLA-DMB, CD86, IL2RA, VAV1, JAG1, JAK3, NOTCH1, HLA-DPA1</i>
Th1 Pathway	6.07	0.131	<i>CCR5, HLA-DOA, CD3E, IL12RB1, HAVCR2, PIK3R5, HLA-DQA1, LGALS9, FGFR2, CD3G, IL18, HLA-DMB, IL10RA, CD86, VAV1, JAK3, NOTCH1, HLA-DPA1</i>
Role of BRCA1 in DNA Damage Response	5.59	0.162	<i>RAD51, FANCB, FANCD2, RFC4, FANCG, BARD1, SMARCD1, PLK1, E2F3, BLM, RBL1, E2F2, CHEK1</i>
CD28 Signaling in T Helper Cells	5.56	0.127	<i>HLA-DOA, ARPC1B, CD3E, PIK3R5, HLA-DQA1, FGFR2, IKBKE, CD3G, LCK, CARD11, SYK, ITPR3, HLA-DMB, FCER1G, CD86, VAV1, LCP2</i>
iCOS-iCOSL Signaling in T Helper Cells	5.32	0.128	<i>HLA-DOA, IL2RG, CD3E, HLA-DQA1, PIK3R5, FGFR2, IKBKE, CD3G, LCK, HLA-DMB, ITPR3, FCER1G, CD86, VAV1, IL2RA, LCP2</i>

Supplementary Table 5. Pathways well-predicted in hepatocellular carcinoma samples. Canonical pathways with the best overlap with the genes for which expression was best predicted (correlation coefficient above 0.4) in liver hepatocellular carcinomas. P-values were calculated using Right-Tailed Fisher's Exact Test.

Ingenuity Canonical Pathways	-log(p-value)	Ratio	Genes
Primary Immunodeficiency Signaling	16.3	0.4	<i>CD19, IL2RG, CD3E, IGLL1/IGLL5, CITA, CD79A, IGHG1, CD8A, TNFRSF13C, CD3D, TAP1, IL7R, LCK, IGHG3, ICOS, ZAP70, IGHM, IGHAI, JAK3, TAP2</i>
Cell Cycle Control of Chromosomal Replication	12.8	0.321	<i>MCM5, MCM6, CDC45, CDT1, CDC6, CDC7, ORC6, CDK1, MCM4, MCM3, PCNA, MCM2, TOP2A, PRIM2, CHEK2, DBF4, MCM7, ORC1</i>
Mitotic Roles of Polo-Like Kinase	10.4	0.258	<i>KIF23, CDC25C, ESPL1, CDC20, PTTG1, PRC1, CDC7, CCNB2, PLK1, CDK1, CCNB1, CDC25B, PLK4, FBXO5, CHEK2, KIF11, CDC25A</i>
Th1 and Th2 Activation Pathway	9.73	0.144	<i>CD247, CD3E, KLRD1, IL12RB1, CXCR3, CD8A, TBX21, IL18R1, IL2RB, RUNX3, IFNG, IL2RG, IKZF1, IL12RB2, CD3D, STAT4, CD3G, LTA, ICOS, GFI1, CXCR6, S1PR1, APH1B, HLA-DOB, IL2RA, PIK3CD, JAK3</i>
Cell Cycle: G2/M DNA Damage Checkpoint Regulation	9.32	0.286	<i>CDC25C, CKS2, YWHAZ, CCNB2, PLK1, AURKA, CDK1, CHEK1, SKP2, CCNB1, CDC25B, CKS1B, TOP2A, CHEK2</i>
Role of CHK Proteins in Cell Cycle Checkpoint Control	8.37	0.246	<i>CDC25C, PLK1, E2F3, CDK1, CHEK1, PCNA, RFC4, E2F1, RFC2, CLSPN, E2F2, CHEK2, E2F8, CDC25A</i>
Estrogen-mediated S-phase Entry	8.24	0.385	<i>CCNA2, CCNE1, E2F1, E2F3, ESR1, E2F8, E2F2, CDK1, SKP2, CDC25A</i>
Th2 Pathway	8.07	0.145	<i>CD247, RUNX3, IFNG, IL2RG, CD3E, IL12RB1, IKZF1, IL12RB2, TBX21, CD3D, STAT4, CD3G, ICOS, GFI1, CXCR6, S1PR1, APH1B, HLA-DOB, IL2RA, PIK3CD, JAK3, IL2RB</i>
Th1 Pathway	7.46	0.146	<i>CD247, RUNX3, IFNG, CD3E, IL12RB1, KLRD1, CXCR3, IL12RB2, TBX21, CD8A, CD3D, IL18R1, STAT4, CD3G, LTA, ICOS, APH1B, HLA-DOB, PIK3CD, JAK3</i>
T Helper Cell Differentiation	6.9	0.192	<i>IL6ST, IFNG, IL2RG, IL12RB1, IL21R, FOXP3, IL12RB2, TBX21, IL18R1, STAT4, ICOS, HLA-DOB, IL2RA, TNFRSF1B</i>

Supplementary Table 6. Pathways well-predicted in breast cancer samples.

Canonical pathways with the best overlap with the genes for which expression was best predicted (correlation coefficient above 0.4) in breast cancer samples. P-values were calculated using Right-Tailed Fisher's Exact Test.

Angio genesis	Hypoxia	DNA repair	Cell cycle	B cells	T cells	Housekeep ing
<i>CDC42</i>	<i>ADM</i>	<i>CETN2</i>	<i>ATM</i>	<i>MASP2</i>	<i>CD79B</i>	<i>RPL32</i>
<i>FGFR1</i>	<i>ALDOA</i>	<i>DDB1</i>	<i>ATR</i>	<i>IGLC7</i>	<i>BTLA</i>	<i>PPIA</i>
<i>FLT1</i>	<i>BHLHE40</i>	<i>DDB2</i>	<i>BUB1</i>	<i>IGLC3</i>	<i>WAS</i>	<i>PGK1</i>
<i>FLT4</i>	<i>BNIP3L</i>	<i>ERCC1</i>	<i>BUB1B</i>	<i>IGKV3-20</i>	<i>CTSC</i>	<i>HMBS</i>
<i>HIF1A</i>	<i>CA12</i>	<i>ERCC2</i>	<i>BUB3</i>	<i>C3</i>	<i>ZAP70</i>	<i>GAPDH</i>
<i>HRAS</i>	<i>CCNG2</i>	<i>ERCC3</i>	<i>CCNA1</i>	<i>IGHG4</i>	<i>FYN</i>	<i>GUSB</i>
<i>ITGAV</i>	<i>CDKN1A</i>	<i>ERCC4</i>	<i>CCNA2</i>	<i>IGKV4-1</i>	<i>ANXA1</i>	<i>TBP</i>
<i>JAG1</i>	<i>CDKN1B</i>	<i>ERCC5</i>	<i>CCNB1</i>	<i>CD74</i>	<i>IFNG</i>	<i>PSMB2</i>
<i>KDR</i>	<i>COL5A1</i>	<i>ERCC8</i>	<i>CCNB2</i>	<i>C8A</i>	<i>PVR</i>	<i>ALB</i>
<i>MAPK14</i>	<i>CP</i>	<i>FEN1</i>	<i>CCND1</i>	<i>IGLL1</i>	<i>C3</i>	<i>HPRT1</i>
<i>NRP1</i>	<i>DDIT3</i>	<i>GTF2H1</i>	<i>CCND2</i>	<i>MLH1</i>	<i>LILRB1</i>	<i>EMC7</i>
<i>NOS3</i>	<i>EDN2</i>	<i>GTF2H3</i>	<i>CCND3</i>	<i>SERPING1</i>	<i>HRAS</i>	<i>RPS27</i>
<i>NFATC4</i>	<i>ENO1</i>	<i>GTF2H5</i>	<i>CCNE1</i>	<i>POU2F2</i>	<i>CD74</i>	<i>RPLP0</i>
<i>PIK3CA</i>	<i>F3</i>	<i>LIG1</i>	<i>CCNE2</i>	<i>LIG4</i>	<i>TRAF6</i>	<i>SDHA</i>
<i>PIK3CB</i>	<i>FOS</i>	<i>MPG</i>	<i>CCNH</i>	<i>HLA-DRB1</i>	<i>LILRB5</i>	<i>ACTB</i>
<i>PIK3CG</i>	<i>GAPDH</i>	<i>PCNA</i>	<i>CDC14A</i>	<i>CFI</i>	<i>HLA-B</i>	AC010970. 1
<i>PIK3R1</i>	<i>HIF1A</i>	<i>POLB</i>	<i>CDC16</i>	<i>IL4R</i>	<i>FCGR1B</i>	
<i>PDGFA</i>	<i>IGFBP3</i>	<i>POLD1</i>	<i>CDC20</i>	<i>CD40LG</i>	<i>PRKCQ</i>	
<i>PRKCA</i>	<i>HK1</i>	<i>POLD3</i>	<i>CDC23</i>	<i>BCL3</i>	<i>HLA-DRB1</i>	
<i>PRKCB</i>	<i>HK2</i>	<i>POLD4</i>	<i>CDC25A</i>	<i>IGLC6</i>	<i>ERAP1</i>	
<i>PLCG1</i>	<i>HMOX1</i>	<i>POLH</i>	<i>CDC25B</i>	<i>RNF8</i>	<i>IL4R</i>	
<i>PTK2</i>	<i>IGFBP1</i>	<i>POLL</i>	<i>CDC25C</i>	<i>IGHM</i>	<i>RIPK2</i>	
<i>PTGS2</i>	<i>IL6</i>	<i>POLR2A</i>	<i>CDC26</i>	<i>HLA-DQB1</i>	<i>IL4</i>	
<i>PXN</i>	<i>JUN</i>	<i>POLR2C</i>	<i>CDC27</i>	<i>C4BPB</i>	<i>CD40LG</i>	
<i>SHC1</i>	<i>LDHA</i>	<i>POLR2D</i>	<i>CDC45</i>	<i>BATF</i>	<i>HMGB1</i>	

<i>SH2D2A</i>	<i>MIF</i>	<i>POLR2E</i>	<i>CDC6</i>	<i>ZP3</i>	<i>KLRK1</i>	
<i>VEGFA</i>	<i>P4HA1</i>	<i>POLR2F</i>	<i>CDC7</i>	<i>IGHG3</i>	<i>CTLA4</i>	
<i>VAV2</i>	<i>PDGFB</i>	<i>POLR2G</i>	<i>CDK1</i>	<i>SWAP70</i>	<i>LILRB2</i>	
	<i>PFKL</i>	<i>POLR2H</i>	<i>CDK2</i>	<i>C1QBP</i>	<i>PTPRC</i>	
	<i>PKP</i>	<i>POLR2I</i>	<i>CDK4</i>	<i>BCL10</i>	<i>HAVCR2</i>	
	<i>PGF</i>	<i>POLR2J</i>	<i>CDK6</i>	<i>FAS</i>	<i>C4BPB</i>	
	<i>PGK1</i>	<i>POLR2K</i>	<i>CDK7</i>	<i>EXO1</i>	<i>SLAMF1</i>	
	<i>PLAUR</i>	<i>RAD51</i>	<i>CDKN1A</i>	<i>IGLC2</i>	<i>ZP3</i>	
	<i>SLC2A1</i>	<i>RAD52</i>	<i>CDKN1B</i>	<i>IGHG2</i>	<i>LILRA1</i>	
	<i>SLC2A3</i>	<i>RBX1</i>	<i>CDKN2A</i>	<i>IGKV2-40</i>	<i>TRAT1</i>	
	<i>STC1</i>	<i>REV3L</i>	<i>CDKN2B</i>	<i>CD55</i>	<i>STAT6</i>	
	<i>TGFB3</i>	<i>RFC2</i>	<i>CDKN2C</i>	<i>PRKCD</i>	<i>PTPN6</i>	
	<i>TGM2</i>	<i>RFC3</i>	<i>CDKN2D</i>	<i>INPP5D</i>	<i>CTSS</i>	
	<i>VEGFA</i>	<i>RFC4</i>	<i>CHEK1</i>	<i>AICDA</i>	<i>BCL10</i>	
		<i>RFC5</i>	<i>CHEK2</i>	<i>IGHV4OR1 5-8</i>	<i>LYN</i>	
		<i>RPA2</i>	<i>CUL1</i>	<i>IL13RA2</i>	<i>TGFB1</i>	
		<i>RPA3</i>	<i>DBF4</i>	<i>C4A</i>	<i>TAP2</i>	
		<i>TP53</i>	<i>E2F1</i>	<i>C1R</i>	<i>TRPM4</i>	
		<i>XPC</i>	<i>E2F2</i>	<i>IGHA2</i>	<i>CD8A</i>	
		<i>ATM</i>	<i>E2F3</i>	<i>C4BPA</i>	<i>SYK</i>	
		<i>ATR</i>	<i>E2F4</i>	<i>NBN</i>	<i>CTSH</i>	
		<i>CHEK1</i>	<i>E2F5</i>	<i>IGHG1</i>	<i>PIK3CD</i>	
		<i>CHEK2</i>	<i>HDAC1</i>	<i>C7</i>	<i>FOXP3</i>	
		<i>H2AFX</i>	<i>KI67</i>	<i>MBL2</i>	<i>CTSL</i>	
			<i>MAD1L1</i>	<i>IGHV2-70</i>	<i>PAG1</i>	
			<i>MAD2L1</i>	<i>C1QA</i>	<i>NECTIN2</i>	
			<i>MCM2</i>	<i>HSPD1</i>	<i>CD79A</i>	

			<i>MCM3</i>	<i>C1RL</i>	<i>C4BPA</i>	
			<i>MCM4</i>	<i>IGKV1D-33</i>	<i>TNFSF18</i>	
			<i>MCM5</i>	<i>C9</i>	<i>CRTAM</i>	
			<i>MCM6</i>	<i>MSH6</i>	<i>BTK</i>	
			<i>MCM7</i>	<i>C8G</i>	<i>IL2</i>	
			<i>MDM2</i>	<i>IGLL5</i>	<i>RAET1E</i>	
			<i>MYC</i>	<i>IGLV7-43</i>	<i>HSPD1</i>	
			<i>ORC1</i>	<i>IRF7</i>	<i>HLA-A</i>	
			<i>PCNA</i>	<i>CD27</i>	<i>MSH6</i>	
			<i>PKMYT1</i>	<i>GCNT3</i>	<i>TNFRSF13C</i>	
			<i>PLK1</i>	<i>C6</i>	<i>EIF2AK4</i>	
			<i>PTTG1</i>	<i>ERCC1</i>	<i>HLA-E</i>	
			<i>RAD21</i>	<i>IGKV3D-20</i>	<i>IRF7</i>	
			<i>RB1</i>	<i>EXOSC3</i>	<i>PRKCB</i>	
			<i>RBL1</i>	<i>FCER1G</i>	<i>CD4</i>	
			<i>RBL2</i>	<i>C1S</i>	<i>EXOSC3</i>	
			<i>SKP1</i>	<i>IGKV1-5</i>	<i>IL12B</i>	
			<i>SKP2</i>	<i>MSH2</i>	<i>FCER1G</i>	
			<i>SMC1A</i>	<i>IGKC</i>	<i>JAK3</i>	
			<i>SMC1B</i>	<i>GAPT</i>	<i>B2M</i>	
			<i>SMC3</i>	<i>C1QC</i>	<i>BCL6</i>	
			<i>STAG1</i>	<i>CR2</i>	<i>LILRB3</i>	
			<i>STAG2</i>	<i>RNF168</i>	<i>MEF2C</i>	
			<i>TFDP1</i>	<i>IGHA1</i>	<i>MALT1</i>	
			<i>TP53</i>	<i>IGHV3-23</i>	<i>LILRB4</i>	
			<i>WEE1</i>	<i>HLA-DRB5</i>	<i>CSK</i>	
				<i>CR1</i>	<i>TAP1</i>	

				<i>IGKV3D-11</i>	<i>HLA-DRB5</i>	
				<i>TRDC</i>	<i>CR1</i>	
				<i>TLR8</i>	<i>GATA3</i>	
				<i>EXOSC6</i>	<i>CD8B</i>	
				<i>C5</i>	<i>CD86</i>	
				<i>C2</i>	<i>SUSD4</i>	
				<i>SLA2</i>	<i>MAP3K7</i>	
				<i>SUSD4</i>	<i>SLC11A1</i>	
				<i>CD46</i>	<i>NLRP10</i>	
				<i>CLU</i>	<i>IFNB1</i>	
				<i>C8B</i>	<i>ORAI1</i>	
				<i>C4B</i>	<i>ITK</i>	
				<i>IGHV1OR2 1-1</i>		
				<i>IGHE</i>		
				<i>C1QB</i>		
				<i>IGHD</i>		
				<i>IGLC1</i>		

Supplementary Table 7. List of genes used to define the six signatures chosen.
 Defining angiogenesis, hypoxia, DNA repair, cell cycle pathways, B-cell and T-cell immune responses.